**Fermilab Petascale Computing and Computation Activities.**
**Workshop on Great Lakes Consortium for Petascale Computation**
**Vicky White, Ruth Pordes, Don Petravick**
**August 18ᵗʰ 2006**

**Applications**
At Fermilab we provide the core computing fabric and contribute to the global computational facilities  in support of the Laboratory's physics research mission. Drivers of the petascale computing and computations needed by this research are: Access to millions of accelerator and detector control and monitoring channels; Archival and management of datasets up to tens of petabytes a year; The complexity and rarity of the physics phenomena resulting in large-scale statistics and $100^{th}$ of a percent level accuracy needs for data processing calculations; Multiple research groups within globally distributed collaborations of thousands working together on common physics activities.

We currently support the Tevatron Run II experiments (CDF, DO) acquiring data at an average of about a Terabyte a day. We are preparing for the LHC experiments, with large US collaborations for ATLAS and CMS that will handle 10s of petabytes per year. We are working to host the International Linear Collider, now in its R&D phase, through significant contributions to the global design and modeling efforts. We also support additional needs for the theoretical physics research of Lattice QCD, and astrophysics and neutrino physics experiments that are part of the Fermilab research program.

CMS, an exemplar application for the future, has the following performance requirements within the next 2-3 years: Round the clock distribution of the data acquired from the detector (225 MB/sec from the online system to the reconstruction farms) quasi real-time from CERN to Fermilab; Storage at Fermilab for more than 30% of the processed, simulation, calibration and analysis data sets, greater than ten Petabytes a year; Data distribution to more than ten university facilities to provide year round regular (~weekly) refresh of  hundreds of terabyes of disk cache, in addition to data distributions between Fermilab and the six regional data repository facilities (Tier-1s); Support for ten thousand jobs per client workflow and fifty thousand jobs a day, with  > 99% success rate for individual job executions.  In addition Fermilab is hosting a remote operations center for the LHC accelerator and experiments (LHC@FNAL), which will require low latency high bandwidth collaboration tools and data access.

These requirements are increased about two fold to support the currently running Tevatron experiments, neutrino and astrophysics experiments, theoretical physics, future accelerator design activities, and the increasing use of simulation and modeling techniques. A larger fraction of these resources are resident at Fermilab than the CMS example.

Today at Fermilab we have thirty thousand square feet of computer rooms, five thousand (soon to be seven thousand) computers, eight automated tape libraries with 4 petabytes of custodial data, 100 Terabytes of distributed disk cache, internal network backbones of up to 200 Gigabit/sec and, by the end of 2006, eight 10Gigabit/sec connections to the external network. We manage more than 10,000 IP connected devices. Data distribution is currently one petabyte a month, with an average annual doubling in the rate.

**Data Management, Worldwide Distribution and Access**
**Archival tape storage**, based on automated robots, continues to provide highest performance to cost solutions for storage and delivery of a hundred of petabytes of data. In the Fermilab Computing Division, we provide operations activities and developments to ensure high

availability of the systems, zero data loss, migration to new tape technologies every five or so years, complete and accurate management of the file, tape and data information, and publication of metrics including I/O rates, tape occupancy and mounts, and data access patterns. Figure 1 shows current information for total data transferred by the tape systems per day for Fermilab users - approximately equally divided in reads and writes, but currently with read predominating by Run II and writes predominating by CMS.
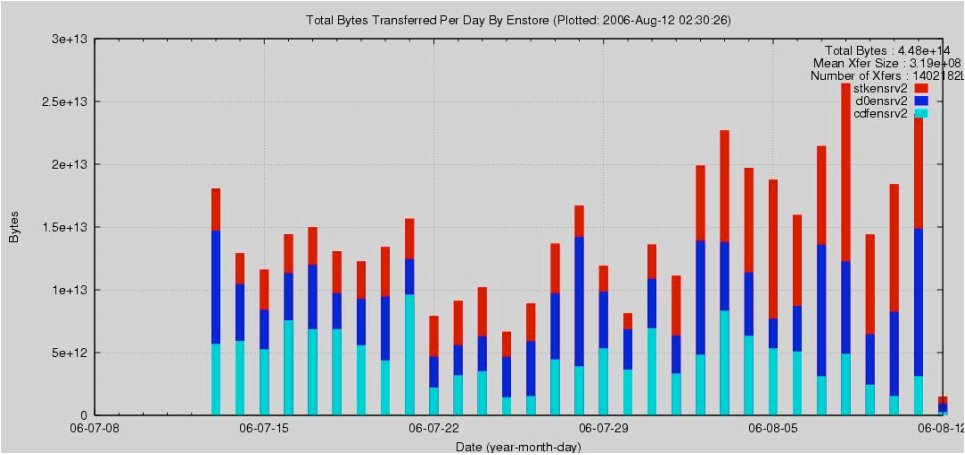


**Figure 1: Enstore MSS system data transfer per day**

**Distributed managed disk caches** of the data needed for data processing and analyses must support bi-directional streaming and random access reads of millions of 1-2 GByte files by hundreds of independent local user jobs. The dCache system (a collaboration between DESY in Germany and Fermilab developers) is providing I/O rates of 200 Mbytes/sec to users for Run II, and up to 200Tbytes a day to the CMS local facility. Figure 2 shows some recent I/O rates of the CMS system. The performance needs to scale up to an order of magnitude to support CMS.
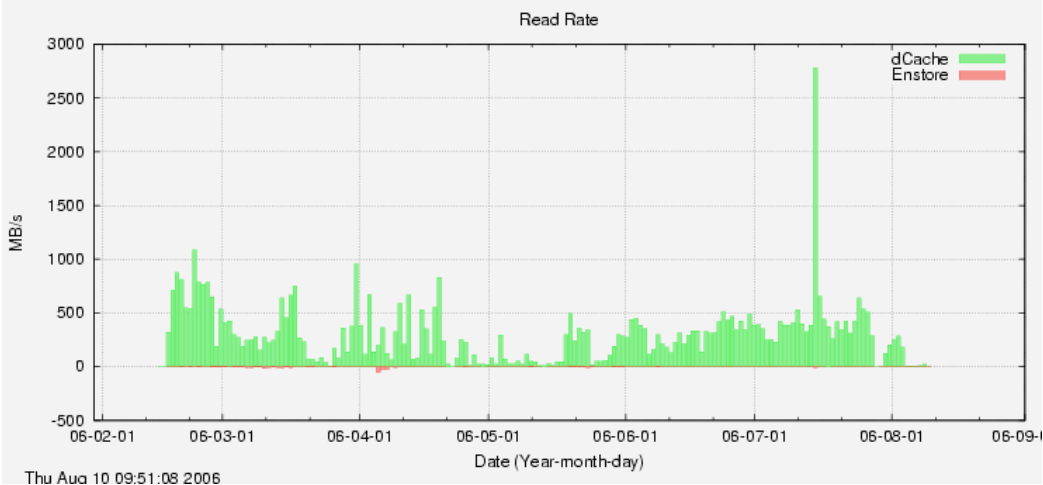


**Figure 2: CMS dCache system I/O rate per day**

Our work at Fermilab includes addressing the following challenges for disk based data caches: Ensure zero data corruption; Maintain quality of service in the face of hardware, software and/or security failures; Manage data movement scheduling, defining and applying policies for data access; Automate file migration and replication to/from the tape backend and across the local and wide area networks; Run operations processes to ensure the highest aggregate throughput.

**Managed networking** is needed to support the wide area distribution of and access to physics data. The LambdaStation software is being developed to provide dynamic allocation of network pipes between available routes (multiple high bandwidth or high/low bandwidth paths). Research and development activities are increasingly focusing on circuit based networks. CMS data challenge activities have shown up needed developments to improve overall stability, response to network glitches, and sustained level of service and device uptimes of >99.99%. We are also working on performance of the end-to-end data transport through closer integration of the network and the storage systems, with an added goal to provide agile allocation of and response to available bandwidth. This is leading to discoveries of unexpected behaviors in the Linux I/O drivers and OS.

For wide area networks, Fermilab uses both ESNET and Starlight and is heavily involved in the management and use of the dedicated LHCnet between CERN and the US as well as the design, implementation and ongoing support for the ESNET Chicago Metropolitan Area Network (MAN). Figure 3 shows increasing offsite network traffic enabled by recent increases in bandwidth capacity, and reaching 1 Petabyte per month for Run II and CMS data distribution.
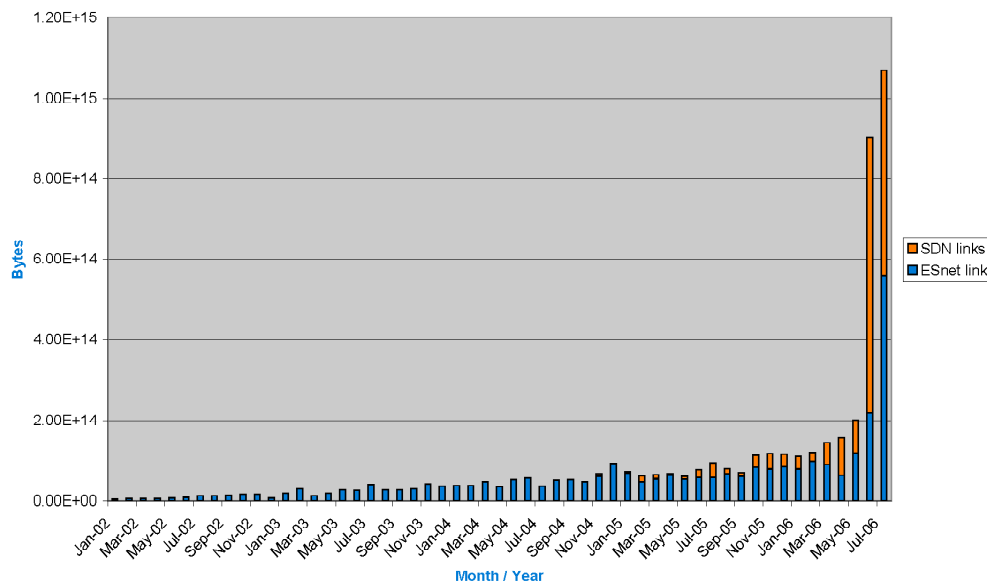


**Figure 3:Fermilab outbound network traffic**

**Computing and Computation**

Fermilab experiment computing and costing models rely on executing jobs on computing farms distributed across the Fermilab facility and the full set of collaborating institutions. We must be cognizant of end-to-end security requirements and incident response of these systems. We must support evolution of both hardware and software over a several years lifecycle. We have addressed these by a strategy that includes well-defined services with well-defined interfaces and scopes as a means have independence of implementation and to facilitate the integration of diverse and evolving components. In synergy with the CERN Worldwide LHC Collaboration, Fermilab has adopted Grids as the means to implement end-to-end distributed systems with its user communities.

Activities in support of the computational services (in addition to the data focused ones) include:
- Comprehensive risk based security management (based on NIST standards), which has been reviewed as one of the "best in class" by the DOE.

- Information management of data location and provenance, job tracking and history, accounting, and performance and operational metrics.
- Facility management infrastructure for hardware and software control, monitoring, fault detection and response.
- Resource management across the aggregated system that reacts to utilization, prioritization and access performance and policies of the specific resources being used.
- Workflow technologies to automate aspects of executing a series of computational steps (e.g. input data placement, job execution, response to error conditions, storage of results).
- Information repositories for inventory, configuration, auditing and tracking of the ensemble of components.
- Distribution and support of Scientific Linux, based on Redhat Enterprise Linux, with common extensions for the scientific community and for distributed multi-site collaborations.

**Integrated Distributed Systems**

The SAMGrid system provides global data management and distribution as well as job management and execution for the Run II experiments. SAMGrid was developed jointly between the Fermilab Computing Division and the experiments. The software system has been maintained and extended over eight years, during experiment commissioning and data taking, to support physics event simulation, processing and analysis. SAMGrid is installed at more than thirty sites worldwide. The heart of SAMGRid is a comprehensive meta-data and information management system, which records the locations and contents of every registered file, and the record of every official experiment job.

Fermilab is integrating its major computing resources into a single logical system (FermiGrid), and is gaining recognition in its management and interfacing of this architecture. FermiGrid allows sharing of, and greater efficiency in, the aggregate use of the farms and storage systems. It also enables coordinated management of job scheduling and execution across the local and remote facilities. Fermilab is also making leadership contributions to the Open Science Grid (OSG), a national distributed infrastructure providing a common set of computational services across existing computing facilities. And, through OSG, Fermilab is an active partner with the Enabling Grids for EsciencE(EGEE) and TeraGrid projects to provide shared cyberinfrastructures based on open standards. Today OSG supports about three thousand simultaneous jobs, across more than twenty-five active facilities, shared by more than eight active research groups.

In addition to these general systems, Fermilab is extending the Lattice QCD facility of tightly coupled clusters, with high bandwidth low latency infiniband connectivity, at a cost of less than $1/Megaflop. In 2005 the 2.2 TF/s Pentium 640 cluster performed in the TOP500. Later in 2006 we will have up to 5 TF/s in production with the addition of dual-core Opterons. This is not a system you can buy on "e-bay" and computations run on the LQCD facility have resulted in significant QCD results acknowledged worldwide.

**References**

1. Fermilab.
2. Fermilab Computing Division.
3. Fermilab Storage Systems.
4. CMS Computing TDR.
5. dCache Collaboration.
6. Computing Division metrics.
7. FermiGrid.
8. Lattice QCD Facility.
9. LambdaStation project.
10. LHC@FNAL
11. Run II Computing Plans 2005.
12. SamGrid.
13. Scientific Linux at Fermilab.
14. Open Science Grid program of work